# New Tools Mobilize Local Data to Study Global Environmental Issues



Youngryel Ryu

**Tracking Carbon Exchange:** Since May 2001, the Flux tower at California's Tonzi Ranch has been tracking rates of carbon dioxide exchange between the atmosphere, plants and soil.

As they strive to develop effective strategies for guarding water supplies, protecting endangered species and curbing greenhouse gases, environmental scientists are turning to innovative cyber-infrastructures and data-mining tools developed by an ongoing collaboration between researchers at Lawrence Berkeley National Laboratory, Microsoft Research, and the University of California, Berkeley.

The Microsoft eScience program is the primary funder of this project, which is one of numerous ventures cultivated by the Berkeley Water Center (BWC). Launched approximately three years ago by researchers from the Berkeley Lab and UC Berkeley's Colleges of Engineering and Natural Resources, the BWC marshals expertise from public institutions and the private sector in support of proj-

ects that enable science and public policy researchers to more easily access and work with water and environmental datasets.

"The most cost-efficient way to impact issues like global climate change and water management is to develop cyber-architectures that organize data and foster scientific collaboration," says Susan Hubbard, staff scientist in the Berkeley Lab's Earth Sciences Division and associate director of the BWC.

Environmental scientists typically collect data on a project-by-project basis, in campaigns targeted at very specific topics. One study may use NASA satellites to track annual rainfall of deserts around the globe, while another project sponsored by the National Science Foundation (NSF) might measure the

# On Tour: CRD Researchers Share Their Expertise at Home and Abroad

CRD researchers are kicking off 2009 by sharing their HPC and networking knowledge and expertise with contemporaries, policy makers and vendors, at home and abroad.

Here's a list of some of the conferences that included CRD staff presentations:

**January 22–23:** Horst Simon, CRD Division Director, appeared on the Greening the Power Hungry Data Center panel at the Greening of the Internet Economy workshop. The workshop was held at the University of California, San Diego campus.

This invitation-only event brought together hundreds of the foremost public policy makers, industry experts, research leaders and global nonprofits to advance strategies for sustainability and energy efficiency in the rapidly growing realm of information and communications technology. The event was hosted by the California Institute for Telecommunications and Information Technology (Calit2) and the California Public Utilities Commission (CPUC).

"Everybody does green. We're halfway there because people have recognized we have a problem," says Simon. "But if we look at vendors and how everybody builds technology, we're not even beginning to make a difference."

For more information on Simon's remarks, visit: http://www.dagstuhl.de/en/program/calendar/semhp/?semnr=09061

**January 26–February 2:** Phil Colella, leader of the SciDAC Applied Partial

## On Tour

Differential Equations Center and head of LBNL's Applied Numerical Algorithms Group, will travel to France to give two talks and meet with researchers at the Commissariat à l'Énergie Atomique (CEA), the French Atomic Energy Commission.

Colella starts his trip at the January 26–27 seminar on numerical fluid mechanics at the Institut Henri Poincaré, where he will give a talk on "Embedded boundary methods and software for solving partial differential equations in complex geometries." After that, he will be meeting with groups at CEA working on modeling and simulation for nuclear reactors. Finally he will give a talk in the Mathematics Department at the University of Paris on February 2.

**January 28, 2009:** Bill Johnston, former department head of the Energy Sciences Network (ESnet), spoke at a live case study hosted by Juniper Networks at the Hyatt Regency Hotel in Reston, Virginia.

In this invitation-only technical seminar for government and commercial Juniper customers, Johnston described the hardware used to develop the architecture and implementation of ESnet4, the organization's next-generation network. This new network is capable of accommodating the massive data flows required for scientific collaborations.

**February 1–6, 2009:** Esmond Ng, Sherri Li, and Horst Simon will participate in the Dagstuhl Seminar on Combinatorial Scientific Computing. Hosted by the Schloss Dagstuhl-Leibniz Center for Informatics, the meeting will be held in Germany. Simon will chair the Monday morning session on February 2, which will include talks on Combinatorial Scientific Computing and Parallel Algorithms in CSC.

Schloss Dagstuhl-Leibniz Center for Informatics promotes fundamental and applied research, continuing and advanced academic education, and the transfer of knowledge between those involved in the research side and application side of informatics.

# The Future of Computer Architecture

*By David Patterson*
**UC Berkeley and CRD's Future Technologies Group**



David Patterson

*(Originally published as part of the Computing Community Consortium's (CCC) Computing Research Initiatives for the 21st Century: http://www.cra.org/ccc/initiatives)*

Computer Architecture is the field that designs computers, which sets the foundation for the entire IT industry.

Despite the tremendous resources at large companies such as IBM and Intel, there is a long track-record of breakthroughs from academic researchers in computer architecture that have led to new multi-billion-dollar industries. One reason is simply that many of the leading thinkers in computer architecture are in academia, as the number of awards, books, and papers document. Another reason is that academic researchers are not bound the business requirements of compatibility with legacy hardware and software. A third reason is that when information technology is changing rapidly, decades of experience may not be a huge asset, and bright young graduate students who don't know what can't be done are as likely to make an important contribution as those considerably more senior.

Thus, in areas where there is tremendous consensus on the guidelines on what to do, computer architects at Intel and IBM usually lead the way. When the directions are less clear, academics are often the path breakers even though they have fewer resources.

Today, we have hit the wall as to the practical limit to the amount of power that a microprocessor chip can dissipate; in the past each generation of chip used more power while getting more performance. We need to invent a new way to get more performance without more power. Moreover, the most interesting computers of the future are at the extremes in scale:

- The datacenter is the computer. Amazon, Google, Microsoft, and others are racing to construct buildings with 50,000 computers to run software as a service (SaaS).
- The cell phone is the computer. Millions of cell phones are shipped each day and they are increasing in functionality.

The "power wall" and the extremes in computer size mean that the old guidelines are out the window, so, if history is any guide, we're entering an era of increasing importance for academic computer architecture.

Note that in such an uncertain era there many chances for new multi-billion-dollar industries, and these new industries are likely to flourish close to where the researchers do their work. We have seen the center of the IT universe shift geographically before during eras of change:

- New York in the mainframe era of the 1950s and 1960s;
- Boston in the minicomputer era of the 1970s and early 1980s;
- Silicon Valley in PC and Web era of the late 1980s to today.

Given the challenges mentioned above, other countries are investing tremendously in IT in general and in computer architecture specifically, since they see the opportunity in this era of change to shift the center of the IT universe and the jobs that come with it. Especially given the cut to academic IT research in general and to computer architecture specifically in the US, it's not hard to imagine the IT center moving further west in the next decade—to Beijing or Mumbai.

One area of tremendous opportunity in computer architecture is the manycore challenge. The goal is to invent computers that make it easy to write programs that are efficient, portable, correct, and

# Environmental Database *continued from page 1*

annual water tables of the Sahara desert with commercial sensors. The data are then typically stored in local archive systems and accessed by researchers associated with that particular project. These sites are scattered across the country, tend to be aligned with specific campaigns, and are funded by a variety of organizations.

According to Catharine van Ingen, partner architect with Microsoft Research, this system can be cumbersome at times because observations are stored in data archives and access centers in the same format that is deposited, and undergo only very simple checks and transformations, making the data difficult to share with other scientists. She notes that much of this information is not science-ready. To fulfill this requirement the data must cataloged, checked, and processed to eliminate obvious problems caused by battery loss, transcription errors, or environmental factors such as freezing rain or birds.

In most cases, scientists also cannot withdraw data from these centers during non-business hours, and so many researchers opt to retain their observations on their own personal desktop computers. If other researchers want to use this data, they have to contact the lead scientist and have him/her e-mail this information to them.

"One of the greatest challenges of the next century will be developing cyber-architectures that allow scientists to easily navigate their digital assets. Today, the internet has given environmental researchers instant access to a wealth of field data. Now, they need a scientific 'safety deposit box' system that will not only store this information, but also organize it so it is searchable and ready for analysis," says van Ingen.

## Designing an Environmental Database for the 21st Century

According to Deb Agarwal, member of the BWC and head of the Advanced Computing for Science Department in the Computational Research Division at Berkeley Lab, the computing needs of many eScience researchers fall into the gap between the typical supercomputer user and the desktop computer user.

"An environmental dataset is often 1 terabyte or smaller in size. These datasets can be stored easily on a desktop hard drive. This means that the hardware needed to create a centralized database is extremely inexpensive and is not the limiting factor. Instead, usability and longevity of the data is the issue," she says.

Agarwal's team worked with existing Microsoft tools initially to develop a prototype database for data collected by the AmeriFlux network. For over 10 years, the AmeriFlux collaboration of field researchers has tracked carbon dioxide exchange between plants and soil on the ground with the planet's atmosphere, on an hourly basis, and in more than 120 sites across North, Central and South America. The sites represent a range of ecosystems, from the Arctic tundra to North American prairies and Amazonian rainforests. Since its inception almost two years ago, the database, called the Fluxdata Scientific Data Server, has grown to include data from Fluxnet, which incorporates AmeriFlux counterparts around the world, including Asia, Africa, Australia, and Europe.

The Fluxdata Scientific Data Server now includes semi-automated ingest tools to automatically extract important aspects of incoming data; a database and schema to organize and archive information; data cubes that allow researchers to



Winter Floods: This aerial photo shows the Wohler ponds along the Russian River during winter as indicated by the high flows and turbid water conditions. The Mirabel inflatable dam is only erected in summer so is not evident in this picture.

*Sonoma County Water Agency*



Watching Water: A water reservoir at Tonzi Ranch in October 2008.

*Youngryel Ryu*

look at the data from multiple perspectives; and tools which automatically convert multiple data versions into one format. The new architecture also enables researchers to browse data and reports via Internet and collaborate with each other. This means scientists no longer need to download and interpret the raw data from a data collection center. Instead they can browse, mine, and do research on the data without needing to download and process it first.

Once this server architecture proved to

## Environmental Database _continued from page 3_

be successful, eScience team members applied this "cyber-blueprint" to create searchable central repositories for the variety of field data collected from California's Russian and Pajaro Rivers. Currently, the team members are collaborating with the National Marine Fisheries Service to aid research involving fish recovery efforts in Northern California coastal streams, and will soon develop a server than encompasses observational information about all the watersheds in California.

"In the past, the computing needs of environmental researchers have often been overlooked because they are rarely on the leading edge of computational or scale requirements of the scientific community, and collectively are not a big enough customer to be commercially profitable. Despite this, their computing challenges are substantial and solving them is essential to their work helping us understand climate change and our surrounding environment," says Agarwal.

CRD currently hosts seven, soon to be eight, BWC Data Servers in the Advanced Computing for Science Department's machine room. The machines are supplied by CRD and Microsoft, while Agarwal's team of CRD scientists provide hardware support, networking, rack space, electricity, and console servers.

Current team members include CRD Scientists Keith Jackson and Monte Goode. Past members of the team include Robin Weber, a UCB employee, and Matt Rodriguez, former CRD member. Kurt Spindler, a high school summer student from the Berkeley Lab's Center for Science & Engineering Education program, also worked on the project in the summer of 2008.

**"One of the greatest challenges of the next century will be developing cyber-architectures that allow scientists to easily navigate their digital assets..."**

**— van Ingen**

### Impacting Science

According to Jim Hunt, professor of civil engineering at UC Berkeley and co-director of the BWC, relatively basic questions such as how the annual water balance in the Russian River watershed changed in the past decade were not exactly impossible to answer before the eScience data-mining tools were developed. However, the tasks of gathering data from a variety of organizations, reformatting the data to make it consistent, sifting out the important pieces of information, and calculating the balances, were so time-consuming and tedious that most scientists didn't want to tackle the issue. He notes that the new eScience tools can produce this answer in minutes. In addition, the data cube architecture allows scientists to find many different relationships in the datasets.

"Everything in an ecosystem is interconnected. Changes in one particular ecosystem could have global consequences, and tools like the data cube make it easier for us to see the big picture.… We can now inquire about more complex relationships like how do the changes in a watershed's annual water balance affect the amount of carbon dioxide in its surrounding atmosphere," says Dennis Baldocchi, Professor of Biometeorology at UC Berkeley.

"The answers to these types of questions will allow us to make accurate predictions about the future of such watersheds, and in turn helps us develop more effective strategies for managing these resources," adds Hunt.

For more information on the BWC, please visit: http://bwc.berkeley.edu/home/thrust_areas/mstci.html

## Architecture _continued from page 2_

scale as the number of cores per microprocessor increases—as easy as it has been to write programs for traditional computers. If this invention allows software to use many simple power-efficient cores instead of a single power-hungry core, this will reset the foundation for the IT industry for at least the next 30 years.

A second opportunity is inventing a new computer architecture that improves computer security and privacy, problems that plague the IT field. Architects could remove many of the vulnerabilities of today's computers if they were not bound by the legacy requirements of compatibility with today's computers. They could also provide new features to make it easier to build fast, secure, low-overhead virtual machines, making it easier and safer for software to migrate between the datacenter and the cell phone.

A third opportunity is invent computers that will remove the performance bottlenecks from new, highly productive programming environments such Ruby or Python. For example, the Ruby on Rails environment allows programmers to invent amazing new computer applications in just 1000 to 2000 lines of code—factors of 10 to 100 less than conventional approaches. Example application areas include personal health care, personal memory assistants, and personal digital educators. Alas, Ruby on Rails performance is factors of 3 to 10 worse than conventional systems. If we can invent computers that allow new programming systems like Ruby or Python to scale up to hundreds of cores while preserving their amazing programming productivity, we could unleash a new round of exciting applications that will lead to new multi-billion-dollar industries, just as we've done so many times in the past.

Intel and Microsoft recently funded two major academic centers to tackle these critical challenges—at UC Berkeley and the University of Illinois. Many other highly credible proposals were received, from strong teams pursuing diverse approaches. A national initiative would fund five more centers of excellence in computer systems and architecture, positioning the United States to maintain its preeminence in the IT field.

## NEWSBYTES:

### New Employee Profile: Sarah Poon

With the availability of cutting-edge science instruments, supercomputers and networks, scientists around the world can collaborate on large-scale experiments that yield a tremendous amount of data. As the new Computer Systems Engineer with the Berkeley Lab's Advanced Computing for Sciences (ACS) group, Sarah Poon will develop applications and visual interfaces to help scientists make sense of this influx of data.

In the next few months, she will dedicate most of her time to developing user-friendly interfaces for the Particle Data Group's (PDG) online version of The Review, which is a compilation and evaluation of the properties of the elementary particles. The compilation is published every other year.

"Scientists are being inundated with so much data that there is really a pressing need for software to help them make sense of this data.… Usability and user interface design are growing areas in computing to support science, and I feel really excited to help raise awareness and advocate for usability at the Lab," she says.

In addition to her work with PDG, Poon will also setting up a website for the 2009 International Conference on Applied Physics, update the SciDAC Outreach Center website, and develop an application to standardize the way bibliographies are presented.

Although Poon is a new hire, she is not new to Berkeley Lab. As a graduate student at the University of California, Berkeley's School of Information, she co-developed a data warehouse and work-flow visualization system with LBNL's Nearby Supernova Factory. She graduated in 2006 and continued her work full time for a year after. She then moved to Southern California, worked in industry, but continued to hear the Bay Area and Berkeley Lab beckoning.

"After I left the lab, I began working with Cecilia Aragon on several papers, based on the work we did at the lab, and I realized how proud I was of my work and how fulfilling it is to work on such interesting problems. And I also really missed living in the Bay area, so I was already looking for an excuse to move back," says Poon.

She now lives in north Berkeley with her husband and two dogs. On her free time, Poon enjoys hiking with her dogs, biking, reading and knitting.

Prior to pursing a master's degree at UC Berkeley in 2004, Poon graduated from UCLA with a bachelor's degree in business economics with a computing specialization. She then spent several years working as a web developer for various E-commerce sites. A year after receiving her undergraduate degree, she began taking computer science classes part time at Cal State Fullerton, eventually completing the core undergrad curriculum. At UC Berkeley, her interests included applied natural language processing, information visualization, and human-computer interaction.

### ESnet Upgrades Science Network

The Energy Sciences Network continues upgrading ESnet4 in the New Year. After spending most of 2008 planning and upgrading its backbone network, ESnet is now turning its attention to its regional networks supporting DOE facilities.

In January 2009, the Bay Area Metropolitan Area Network (BAMAN) received the most attention as the deployed Cisco 6509 switchers and routers were upgraded to the latest models from Juniper's MX series of routers. The Berkeley Lab and ESnet's Qwest Sunnyvale hub received new Juniper MX960 routers. Meanwhile the Sandia National Laboratory, Lawrence Livermore National Laboratory and Stanford Linear Accelerator Center received new MX480 routers. Testing and configuration of the same model is currently underway for the hardware upgrade at the Joint Genome Institute (JGI), and the NERSC Oakland Scientific Facility will get a MX960 in February.

Approximately 1,800 miles away from California's Bay Area at the Fermi National Laboratory (FNAL) in Bativa, Ill., ESnet will also replace its current Cisco hardware with new Juniper MX960 and MX480 routers in February.

These upgrades provide a foundation for additional optical 10-gigabit interface circuits to be installed in the future. Each 10-gigabit line is capable of transferring the equivalent of 500-hours of digital music per second, which is essential for connecting DOE researchers to future large science experiments like the Large Hadron Collider (LHC).

January 2009 also brought circuit upgrades to ESnet's Denver hub, where a peering connection to the Front Range GigaPop (FRGP) was upgraded from 1-gigabit to 10-gigabit. The FRGP is a network that connects a consortium of universities, non-profit organizations and government agencies, in Colorado and neighboring states. The peering connection links FRGP users to DOE sponsored science transferred over ESnet, and the recent circuit upgrades give users on both networks more bandwidth to exchange data.

ESnet engineers are also currently working with Lightower to get two additional 10-gigabit waves between ESnet hubs in New York City to the Brookhaven National Laboratory (BNL).

## About CRD Report

CRD Report, which publishes every other month, highlights the cutting-edge research conducted by staff scientists in areas including turbulent combustion, nanomaterials, climate change, distributed computing, high-speed networks, astrophysics, biological data management and visualization. CRD Report Editor Linda Vu can be reached at 510 495-2402 or LVu@lbl.gov. Find previous CRD Report articles at http://crd.lbl.gov/html/news/CRDreport.html.